

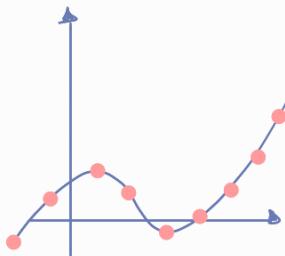
Kernelized Linear Regression

We define a set of functions ϕ_i $i=0, \dots, n$ (in total there are $n+1$) and consider such model:

$$y(x; w) = w^T \phi(x) \quad \text{with } \phi(x) = \begin{bmatrix} \phi_0(x) \\ \vdots \\ \phi_n(x) \end{bmatrix} \quad \text{and } \phi_0(x) = 1 \quad \text{still } n \text{ degrees of freedom}$$

The resulting model is still linear in w but not in x (ϕ is usually non linear, otherwise this whole thing will be pointless).

e.g. we consider this particular transformation: $\phi = [x^0, x^1, \dots, x^n]^T$
 $y(x; w) = w^T \phi(x) = (w_0, w_1, \dots, w_n)(x^0, x^1, \dots, x^n)^T = w_0 + w_1 x + w_2 x^2 + \dots + w_n x^n$
 this is an M -degree polynomial



Solution of the problem is to find a good M such that we will best fit this data
 In this case we could choose $M=3$

The choice of M is critical, it could lead to OVERFITTING:
 higher M could better fit the available data but we are interested in classifying new samples

Target values t in the dataset are given by the true function $y(x; w)$ affected by additive noise ϵ .

$$t = y(x; w) + \epsilon$$

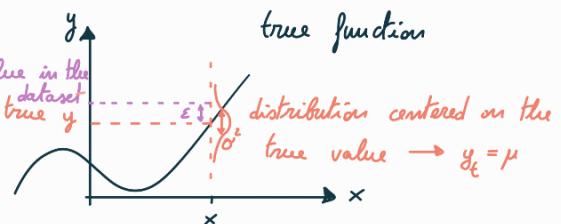
We assume that the noise follows a gaussian distribution.

We just described the same problem from a probabilistic point of view.

$$P(t | x, w, \beta) = N(\epsilon | y(x; w), \beta^{-1}) \quad \text{variance is described in terms of precision } \beta;$$

likelihood

the distribution is centered on the true value (zero-mean, we can translate it (?)



Now we can try to maximize the likelihood (minimize the negative log likelihood)

$$\text{maximize } P(\{t_1, \dots, t_N\} | x_1, \dots, x_N, w, \beta) = \prod_{n=1}^N N(t_n | w^T \phi(x_n), \beta^{-1}) \quad \text{under the hypothesis that the samples are independent and identically distributed we can use the theorem of total probability}$$

or equivalently:

$$\begin{aligned} \text{minimize } \ln P(\{t_1, \dots, t_N\} | x_1, \dots, x_N, w, \beta) &= \sum_{n=1}^N \ln N(t_n | w^T \phi(x_n), \beta^{-1}) \\ &= -\beta \frac{1}{2} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2 - \frac{N}{2} \ln (2\pi \beta^{-1}) \end{aligned}$$

$E_D(w)$
do not depend on w

$E_D(w)$ is what we need to minimize

Maximum likelihood (minimum log likelihood) under zero-mean Gaussian assumption correspond to least square error minimization

$$\underset{w}{\operatorname{argmax}} P(\{t_1, \dots, t_N\} | x_1, \dots, x_N, w, \beta) \longleftrightarrow \underset{w}{\operatorname{argmin}} E_D(w) = \underset{w}{\operatorname{argmin}} \frac{1}{2} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2$$

sum of the squared errors:

sum over all the points in the dataset of the squared difference between the true value t_n and the prediction $w^T \phi(x_n)$

linear regression can be solved with this optimization function that derives at the same time from the probabilistic formulation of the problem (maximum likelihood) and from the least square minimisation (squared difference between the true value t_n and the prediction $w^T \phi(x_n)$).

From the closed form solution:

$$\textcircled{1} \text{ compute the gradient } E_D(w) = \frac{1}{2} (t - \phi w)^T (t - \phi w)$$

with $t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix}$, $\phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_n(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_n(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \dots & \phi_n(x_N) \end{bmatrix}$, $w = \begin{bmatrix} w_0 \\ \vdots \\ w_M \end{bmatrix}$

$$\textcircled{2} \quad \nabla E_D(w) = 0 \iff \phi^T \phi w = \phi^T t$$

$$\textcircled{3} \quad \text{isolate } w \text{ from } \phi^T \phi w = \phi^T t$$

$$w = (\phi^T \phi)^{-1} \phi^T t$$

pseudo-inverse of $\phi \hat{=} \phi^+$

example

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$D = \{(x_m, t_m)\}_{m=1}^N$$

define t, w and ϕ

we assume this particular transformation:

$$\phi = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ & \vdots & & \\ N-1 & x_N & x_N^2 & x_N^3 \\ \uparrow & \uparrow & \uparrow & \uparrow \\ 0 & x_0 & x_0^2 & x_0^3 \end{bmatrix}$$

the ϕ matrix will have N rows, one for each sample in the dataset and $M+1$ cols (number of features of the ϕ vector + ϕ_0)

$$t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix}$$

$$w = \begin{bmatrix} w_0 \\ \vdots \\ w_M \end{bmatrix}$$

our model will be:

$$y(x; w) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

$$(W_{ML}) = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{pmatrix} = \phi^+ t \quad \text{pinv}(\phi) \cdot t$$

w from maximum likelihood

We need an iterative algorithm. Why?

Because the ϕ matrix may be too large (N depends on the dataset size and M could also be very large)

Iterative algorithm: sequential learning

$$\hat{w} \leftarrow \hat{w} - p \nabla E_m$$

$$\hat{w} \leftarrow \hat{w} - p(t_m - w^T \phi(x_m)) \phi(x_m)$$

this algorithm converges for suitable small values of p ; its value need to be chosen carefully

Regression problem

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$D = \{(x_m, t_m)\}_{m=1}^N$$

for linear regression as well as for all other linear models, the functions we end up using for classification are linear in the weight but could (and should) be non-linear in the input

Linear regression

$$y(x; \omega) = \omega_0 + \omega^T x \quad x, \omega \in \mathbb{R}^d$$

Linear regression with non-linear basis function

$$y(x; \omega) = \omega^T \phi(x)$$

$$\phi = [\phi_0(x), \phi_1(x), \dots, \phi_n(x)]^T, \phi_0(x) = 1$$

most commonly used basis function

polynomial
radial
sigmoid

Kernelized linear regression

linear regression is too simple for many applications
Start from a polynomial linear regressor:

$$y = \omega^T \phi(x)$$

$$\phi(x), \omega \in \mathbb{R}^n$$

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^n$$

polynomial regression
f

non-linear regression

polynomial \rightleftarrows non-linear

$$y(x; \omega) = \omega^T \phi(x) = (\omega_0, \omega_1, \dots, \omega_n)(x^0, x^1, \dots, x^n)^T$$

$$\phi_n(x) = x^n$$

How do we find the best ω ?

Define a cost function:

$$J(\omega) = \sum_{m=1}^N (y_m - \omega^T x_m)^2 \text{ least squares cost function}$$

Minimize the cost function with respect to the weights

$$\omega^* = \underset{\omega}{\operatorname{argmin}} J(\omega)$$

After minimizing and simplifying we get:

$$\omega = (X^T X)^{-1} X^T Y$$

For the non-linear basis function:

$$\omega = (\phi^T \phi)^{-1} \phi^T Y$$

With regularization:

$$\omega = (\phi^T \phi + \lambda I)^{-1} \phi^T Y$$

with regularization we talk about
Ridge Regression

Estimating the value of w in terms of ϕ is very computationally expensive ($\phi^T \phi$) but there is another way (**Kernelization**)

lets start again with the Ridge Regression cost function

$$J(w) = \sum_{m=1}^N (y_m - w^T \phi(x_m))^2 + \frac{1}{2} \lambda \sum_{m=1}^N \|w\|^2$$

Take the derivative with respect to the weight vector and equate it to zero

$$\begin{aligned} w^* &= \left(\frac{1}{\lambda} \right) \sum_{m=1}^N (y_m - w^T \phi(x_m)) \phi(x_m) \\ w^* &= \sum_{m=1}^N \boxed{\frac{1}{\lambda} (y_m - w^T \phi(x_m)) \phi(x_m)} \\ &\quad \downarrow \alpha_m \\ w^* &= \sum_{m=1}^N \alpha_m \phi(x_m) \\ &= \phi^T \alpha \quad \text{dot product notation} \end{aligned}$$

Go back to the original cost function and vectorize w

$$J(w) = \sum_{m=1}^N (y_m - w^T \phi(x_m))^2 + \frac{1}{2} \lambda \sum_{m=1}^N \|w\|^2$$

$$J(w) = (y - \phi w)^T (y - \phi w) + \frac{\lambda}{2} w^T w$$

Simplify

$$\begin{aligned} J(w) &= y^T y - y^T \phi w - w^T \phi^T y + w^T \phi^T \phi w + \frac{\lambda}{2} w^T w \\ J(w) &= y^T y - \underbrace{y^T \phi (\phi^T \alpha)}_{\text{cancel}} - \underbrace{(\phi^T \alpha)^T \phi^T y}_{\text{cancel}} + \underbrace{(\phi^T \alpha)^T \phi^T \phi (\phi^T \alpha)}_{\text{cancel}} + \frac{\lambda}{2} (\phi^T \alpha)^T (\phi^T \alpha) \quad w = \phi^T \alpha \end{aligned}$$

$$J(w) = y^T y - \cancel{y^T \phi \phi^T \alpha} - \cancel{\alpha^T \phi \phi^T y} + \cancel{\alpha^T \phi \phi^T \phi \phi^T \alpha} + \frac{\lambda}{2} \cancel{\alpha^T \phi \phi^T \alpha}$$

$\phi \phi^T$ is called Gram matrix / Kernel matrix (K)

$$\phi \phi^T = K = \begin{bmatrix} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \cdots & \phi(x_1)^T \phi(x_m) \\ \phi(x_2)^T \phi(x_1) & \phi(x_2)^T \phi(x_2) & \cdots & \phi(x_2)^T \phi(x_m) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(x_m)^T \phi(x_1) & \phi(x_m)^T \phi(x_2) & \cdots & \phi(x_m)^T \phi(x_m) \end{bmatrix}$$

This matrix is significant because of its two properties :

① Symmetry : $K = K^T$

② Positive semi-definite : $\alpha^T K \alpha \geq 0$ (product with any other vector and its transpose is ≥ 0)

Let's use these properties in our simplification of the cost function

$$J(w) = y^T y - \cancel{y^T \phi \phi^T \alpha} - \cancel{\alpha^T \phi \phi^T y} + \cancel{\alpha^T \phi \phi^T \phi \phi^T \alpha} + \frac{\lambda}{2} \cancel{\alpha^T \phi \phi^T \alpha}$$

$$J(w) = y^T y - y^T K \alpha - \cancel{\alpha^T K y} + \alpha^T K^2 \alpha + \frac{\lambda}{2} \alpha^T K \alpha \quad \because K = K^T$$

$$J(w) = y^T y - y^T K \alpha - \cancel{y^T K \alpha} + \alpha^T K^2 \alpha + \frac{\lambda}{2} \alpha^T K \alpha \quad \text{Transpose the scalar terms}$$

$$J(w) = y^T y - 2y^T K \alpha + \alpha^T K^2 \alpha + \frac{\lambda}{2} \alpha^T K \alpha$$

$$\begin{aligned}
\frac{\partial J(\alpha)}{\partial \alpha} &= 0 - 2y^T K + \alpha^T K^2 + \frac{\lambda}{2} \alpha^T K \\
-2y^T K + \alpha^T K^2 + \frac{\lambda}{2} \alpha^T K &= 0 \\
K(-2y^T + \alpha^T K + \frac{\lambda}{2} \alpha^T) &= 0 \\
-2y^T + \alpha^T K + \frac{\lambda}{2} \alpha^T &= 0 \\
\alpha^T K + \frac{\lambda}{2} \alpha^T &= 2y^T \\
\alpha^T (K + \frac{\lambda}{2} I) &= 2y^T \\
\alpha^T &= 2y^T (K + \frac{\lambda}{2} I)^{-1} \\
\alpha &= (K + \frac{\lambda}{2} I)^{-1} 2y \quad \text{we can drop the scaling factor 2} \\
\alpha^* &= (K + \lambda' I)^{-1} y
\end{aligned}$$

We previously had:

$$w = \phi^T \alpha \longrightarrow w = \phi^T (K + \lambda' I)^{-1} y$$

Not kernelized: $w = (\phi^T \phi + \lambda I)^{-1} \phi^T y$

Kernelized: $w = \phi^T (\phi \phi^T + \lambda' I)^{-1} y$

$\phi \phi^T$ is easier to compute than $\phi^T \phi$

According to Mercer's theorem, a symmetric function $K(x, y)$ can be expressed as an inner product $K(x, y) = \langle \phi(x), \phi(y) \rangle$ for some ϕ if and only if $K(x, y)$ is positive semidefinite, i.e.

$$\int K(x, y) g(x) g(y) dx dy > 0 \quad \forall g$$

or equivalently:

$$\begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots \\ K(x_2, x_1) & \ddots & \vdots \\ \vdots & \vdots & \ddots \end{bmatrix} \text{ is positive semidefinite for any collection } \{x_1, \dots, x_n\}$$

Thanks to Mercer's theorem, we can conclude that we can rewrite every term in the Gram matrix K as a function $K(x_i, x_j)$

$$\begin{array}{l} \textcircled{1} \text{ Symmetry: } K = K^T \\ \textcircled{2} \text{ Positive semi-definite: } \alpha^T K \alpha \geq 0 \end{array} \longleftrightarrow K(x, y) = \langle \phi(x), \phi(y) \rangle$$

we had the inner product, we switch to the function instead

$$\phi^T \phi = K = \begin{bmatrix} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \dots & \phi(x_1)^T \phi(x_m) \\ \phi(x_2)^T \phi(x_1) & \phi(x_2)^T \phi(x_2) & \dots & \phi(x_2)^T \phi(x_m) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(x_m)^T \phi(x_1) & \phi(x_m)^T \phi(x_2) & \dots & \phi(x_m)^T \phi(x_m) \end{bmatrix} = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_m) \\ K(x_2, x_1) & K(x_2, x_2) & \dots & K(x_2, x_m) \\ \vdots & \vdots & \ddots & \vdots \\ K(x_m, x_1) & K(x_m, x_2) & \dots & K(x_m, x_m) \end{bmatrix}$$

Consider this polynomial non-linear basis function

$$\phi : x \longrightarrow \phi(x) = \begin{bmatrix} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{bmatrix}$$

We have a kernel matrix where every term is an inner product of two samples

$$\begin{aligned} \phi(x_m)^T \phi(x_n) &= \begin{bmatrix} x_{m1}^2 \\ \sqrt{2} x_{m1} x_{m2} \\ x_{m2}^2 \end{bmatrix}^T \begin{bmatrix} x_{n1}^2 \\ \sqrt{2} x_{n1} x_{n2} \\ x_{n2}^2 \end{bmatrix} \\ &= x_{m1}^2 x_{n1}^2 + 2 x_{m1} x_{m2} x_{n1} x_{n2} + x_{m2}^2 x_{n2}^2 \\ &= (x_{m1} x_{n1} + x_{m2} x_{n2})^2 \\ &= (x_m^T x_n)^2 \\ &= k(x_m, x_n) \end{aligned}$$

This in particular is an example of a polynomial kernel function

We can compute the kernel matrix K without knowing the true nature of ϕ , only knowing the Kernel function k

Expanding the product we can see that it can be represented as a function of the basis vector inputs x_m and x_n

How can we make predictions?

$$y(x) = w^T \phi(x)$$

$$= y (\phi \phi^T + \lambda' I)^{-1} \phi^T \phi(x)$$

y in the dataset

we can express $\phi^T \phi(x)$ as a column vector of the kernel matrix

$$\phi^T \phi(x) = \begin{bmatrix} \phi(x_1)^T \phi(x) \\ \phi(x_2)^T \phi(x) \\ \vdots \\ \phi(x_n)^T \phi(x) \end{bmatrix} = k_x$$

$$= y (\mathbf{k} + \lambda' \mathbf{I})^{-1} k_x \text{ independent of } \phi$$

Even with complex basis functions we can make predictions with only the base features